

Approximate Distance Oracles with Improved Preprocessing Time

Christian Wulff-Nilsen *

Abstract

Given an undirected graph G with m edges, n vertices, and non-negative edge weights, and given an integer $k \geq 1$, we show that for some universal constant c , a $(2k-1)$ -approximate distance oracle for G of size $O(kn^{1+1/k})$ can be constructed in $O(\sqrt{k}m + kn^{1+c/\sqrt{k}})$ time and can answer queries in $O(k)$ time. We also give an oracle which is faster for smaller k . Our results break the quadratic preprocessing time bound of Baswana and Kavitha for all $k \geq 6$ and improve the $O(kmn^{1/k})$ time bound of Thorup and Zwick except for very sparse graphs and small k . When $m = \Omega(n^{1+c/\sqrt{k}})$ and $k = O(1)$, our oracle is optimal w.r.t. both stretch, size, preprocessing time, and query time, assuming a widely believed girth conjecture by Erdős.

*Department of Mathematics and Computer Science, University of Southern Denmark, koolooz@diku.dk, <http://www.imada.sdu.dk/~cwn/>.

1 Introduction

Computing shortest path distances in graphs is a fundamental algorithmic problem and has received a lot of attention for several decades. Classical algorithms include that of Dijkstra which handles graphs with non-negative edge weights, and the algorithm of Bellman-Ford which is slower but applies also when negative edge weights are present. A considerable drawback of these algorithms is that they are too slow for many applications. For instance, a GPS system needs to report shortest path distances extremely fast in very large road networks. If Dijkstra's algorithm is used, in the worst case the entire network would need to be explored just to compute a single distance. Another problem is that the whole graph would have to be stored in memory. If the graph is dense, it might not fit in the main memory and this would slow down computations considerably.

One way to speed up computations is to precompute distances between all pairs of vertices in a preprocessing step and store them all in a look-up table. Distance queries can then be answered in constant time. The fastest known all-pairs shortest paths algorithm has only marginally subcubic running time [5]. For sparser graphs, repeated applications of Dijkstra's algorithm yield $O(mn + n^2 \log n)$ time. Pettie [12] gave a slightly improved bound of $O(mn + n^2 \log \log n)$. Even for sparse graphs, these algorithms are too slow for many applications. Another disadvantage of this scheme is the huge amount of memory required to store the look-up table for all the distances.

1.1 Approximate distances

A way to deal with these issues is to settle for *approximate* shortest path distances. For a directed graph G , a distance from a vertex u to a vertex v along some path in G is of *stretch* $\delta \geq 1$ if the path is at most δ times longer than a shortest path from u to v in G .

Zwick [17] showed how to compute all-pairs stretch $(1 + \epsilon)$ -distances in directed graphs in $\tilde{O}(n^{2.376})$ time for an arbitrarily small constant $\epsilon > 0$. In the seminal paper of Thorup and Zwick [15], it was shown how to preprocess an undirected graph in close to linear time to build a data structure of near-linear size which reports small-stretch distances very fast. More precisely, for an undirected graph G with non-negative edge weights, m edges, and n vertices and for any integer $k \geq 1$, a data structure of size $O(kn^{1+1/k})$ can be built in $O(kmn^{1/k})$ time which gives distances of stretch at most $2k - 1$ in $O(k)$ time. Since this data structure has constant query time for $k = O(1)$, it is referred to as an *approximate distance oracle*. We emphasize that the result only holds for undirected graphs; as shown by Thorup and Zwick, no compact distance oracles exist in general for directed graphs.

The Thorup-Zwick oracle is randomized. Roditty, Thorup, and Zwick [13] showed how to obtain a deterministic algorithm with only a polylogarithmic increase in preprocessing time. Mendel and Naor [10] showed how to improve query time of the Thorup-Zwick oracle to $O(1)$ and space to $O(n^{1+1/k})$ at the cost of a larger stretch of $O(k)$ and a longer preprocessing time of $O(n^{2+1/k} \log n)$.

For small k , the size/stretch tradeoff of the Thorup-Zwick oracle is essentially optimal due to a (widely believed and partially proved) girth conjecture of Erdős from 1963 [9]. The only possible improvement is thus in the time for preprocessing. It was shown by Baswana and Kavitha [2] that an oracle with the same size and stretch can be computed in $O(n^2)$ time which is an improvement when the number of edges m is $\omega(n^{2-1/k}/k)$ and it is the first essentially optimal oracle for $m = \Theta(n^2)$.

1.2 Contributions of this paper

The main contribution of this paper is to break the quadratic time bound of Baswana and Kavitha [2] for $k \geq 6$ and $m = o(n^2)$ and to break the Thorup-Zwick bound [15] when the graph is not too sparse and k not too small. We show that there exists a constant c such that for any integer $k \geq 1$, a $(2k - 1)$ -approximate distance oracle of size $O(kn^{1+1/k})$ and with query time $O(k)$ can be constructed in $O(\sqrt{k}m + kn^{1+c/\sqrt{k}})$ time. When $m = \Omega(n^{1+c/\sqrt{k}})$ and $k = O(1)$, preprocessing time is linear and our construction is thus optimal in *every* respect (stretch, size, preprocessing time, and query time), assuming the girth conjecture. The oracle of Baswana and Kavitha only has linear preprocessing for $m = \Theta(n^2)$ and the $O(kmn^{1/k})$ preprocessing of Thorup and Zwick is super-linear for any constant k .

We also present an oracle which is faster for smaller k . When $k \geq 3$, $k \bmod 3 = 0$, its preprocessing time is $O(km + kn^{3/2+2/k})$. We get similar bounds for other values of $k \geq 3$: when $k \bmod 3 = 1$, preprocessing is $O(km + n^{3/2+1/(2k)+3/(2(k-1))})$ and when $k \bmod 3 = 2$, it is $O(km + n^{3/2+2/(k-2)-1/(k(k-2))})$. In particular, we achieve subquadratic preprocessing for all $k \geq 6$.

Our algorithms are very simple to describe and analyze, given previous black boxes. As in earlier approaches, we make use of random vertex sampling. We apply a result of Baswana and Kavitha [2] to sparsify our graph w.r.t. this sampling while preserving *exact* distances between pairs of vertices that are close in some sense. We construct a Thorup-Zwick oracle for this sparser subgraph, allowing us to report $(2k - 1)$ -approximate distances between such vertex pairs. For pairs that are farther apart, we make use of spanners. For a $\delta \geq 1$, a δ -*spanner* of a graph G is a subgraph that spans all vertices and preserves distances between all vertex pairs up to a factor of δ . To construct our oracle with small preprocessing time for small k , we run the linear-time algorithm of Baswana and Sen [4] to get a spanner of small stretch. We apply Dijkstra to get exact distances in this spanner between all pairs of sampled vertices. These distances together with distances from vertices to their nearest sampled vertex in the original graph allows us to report stretch $(2k - 1)$ -distances also for vertices far apart.

Our oracle with near-linear preprocessing time for larger k does not run Dijkstra but instead constructs, on top of a spanner, what we call a *restricted* oracle. This oracle only allows us to query distances between sampled vertices but is more compact than the Thorup-Zwick oracle. We pick the stretch of the spanner and the restricted oracle to be $\Theta(\sqrt{k})$ and we show that the oracle gives stretch $(2k - 1)$ -distances in the underlying graph.

We have summarized previous results on distance oracles as well as our new results in Table 1.

1.3 Related work

A problem related to distance oracles is that of finding spanners. We have already mentioned the linear-time algorithm of Baswana and Sen [4] to find a spanner of stretch $2k - 1$. There has also been interest in so-called (α, β) -spanners, where α and β are real numbers. Such a spanner H of a graph G ensures that for all vertices u and v , $d_H(u, v) \leq \alpha d_G(u, v) + \beta$. In other words, H allows an *additive* stretch in addition to a multiplicative stretch. Thorup and Zwick [16] showed the existence of $(1 + \epsilon, \beta)$ -spanners of size $O(n^{1+1/k})$ for any constant k , where $\beta = (c/\epsilon)^k$ for some constant c . A $(1, 2)$ -spanner of size $\tilde{O}(n^{3/2})$ was presented by Dor, Halperin, and Zwick [7]. The size was later improved slightly by Elkin and Peleg to $O(n^{3/2})$ [8]. Baswana, Kavitha, Mehlhorn, and Pettie [3] gave a spanner of size $O(n^{4/3})$.

Table 1: Performance of distance oracles in weighted undirected graphs.

Stretch	Query time	Space	Preprocessing time	Reference
1	$O(1)$	$O(n^2)$	$O(mn + n^2 \log \log n)$	[12]
$1 + \epsilon$	$O(1)$	$O(n^2)$	$\tilde{O}(n^{2.376})$	[17]
2	$O(1)$	$O(n^2)$	$O(n^{3/2} \sqrt{m} \log n)$ $O((m\sqrt{n} + n^2) \log n)$	[6] [2]
	$O(1)$	$O(n^{5/3})$	–	[11]
$7/3$	$O(1)$	$O(n^2)$	$O(n^{7/3} \log n)$	[6]
			$O((m^{2/3}n + n^2) \log n)$	[2]
3	$O(1)$	$O(n^2)$	$O(n^2 \log n)$	[6]
	$O(1)$	$O(n^{3/2})$	$O(m\sqrt{n})$	[15]
	$O(k)$	$O(n^{3/2})$	$O(\min\{m\sqrt{n}, kn^{2+\frac{1}{2k}}\})$	[2]
	$O(\sqrt{m})$	$O(m + n)$	–	[11]
$2k - 1$	$O(k)$	$O(kn^{1+\frac{1}{k}})$	$O(kmn^{1+\frac{1}{k}})$ $O(\sqrt{k}m + kn^{1+\frac{c}{\sqrt{k}}})$	[15] this paper
$2k - 1$ ($k \geq 3$)	$O(k)$	$O(kn^{1+\frac{1}{k}})$	$O(\min\{n^2, kmn^{1+\frac{1}{k}}\})$	[2]
$2k - 1$ ($k \geq 3, k \bmod 3 = 0$)	$O(k)$	$O(kn^{1+\frac{1}{k}})$	$O(km + kn^{\frac{3}{2}+\frac{2}{k}})$	this paper
$2k - 1$ ($k \geq 3, k \bmod 3 = 1$)	$O(k)$	$O(kn^{1+\frac{1}{k}})$	$O(km + kn^{\frac{3}{2}+\frac{1}{2k}+\frac{3}{2(k-1)}})$	this paper
$2k - 1$ ($k \geq 3, k \bmod 3 = 2$)	$O(k)$	$O(kn^{1+\frac{1}{k}})$	$O(km + kn^{\frac{3}{2}+\frac{2}{k-2}-\frac{1}{k(k-2)}})$	this paper
$O(k)$	$O(1)$	$O(n^{1+1/k})$	$O(n^{2+1/k} \log n)$	[10]
	$O(k)$	$O(kn^{1+1/k})$	$O(km + kn^{1+1/k+\epsilon})$	this paper

which has additive stretch 6 and no multiplicative stretch. This is currently the smallest known spanner with constant additive stretch and no multiplicative stretch.

Allowing additive stretch in oracles has also been considered. For unweighted graphs, Baswana, Gaur, Sen, and Upadhyay [1] showed how to get subquadratic construction time by allowing constant additive stretch in addition to a multiplicative stretch of k . Pătraşcu and Roditty [11] showed that for unweighted graphs there exists an oracle of size $O(n^{5/3})$ which has multiplicative stretch 2 and additive stretch 1. Furthermore, they showed that for weighted graphs with $m = n^2/\alpha$ edges, there is an oracle of size $O(n^2/\alpha^{1/3})$ with multiplicative stretch 2 and no additive stretch. Preprocessing time was not considered in [11].

The organization of the paper is as follows. In Section 2, we give some basic definitions and notation as well as some tools that will prove useful. In Section 3, we give the stretch $(2k - 1)$ -oracle which is fast for small k . Our near-linear time oracle is then presented in Section 4. Finally, we make some concluding remarks in Section 5.

2 Definitions, Notation, and Toolbox

Let $G = (V, E)$ be an undirected graph with non-negative edge weights. For our problem, we may assume that all edges have strictly positive weight since zero-weight edges can always be contracted. Also, we shall only consider connected graphs. For vertices $u, v \in V$, we denote by $d_G(u, v)$ the shortest path distance in G between u and v .

For a real value $\delta \geq 1$, a δ -spanner of G is a subgraph $H = (V, E_H)$ of G spanning all its vertices such that for any distinct vertices u and v , $d_H(u, v) \leq \delta d_G(u, v)$.

Let S be a non-empty subset of V . For a vertex u , let $p_S(u)$ be the vertex of S nearest to u w.r.t. d_G (ties are resolved arbitrarily). We denote by $B_S(u)$ the set of vertices v with $d_G(u, v) < d_G(u, p_S(u))$. Let $E_S(v)$ denote the set of edges incident to v with weight less than $d_G(v, p_S(v))$. We define $E_S = \cup_{v \in V} E_S(v)$ and $G_S = (V, E_S)$. We need the following two simple results.

Lemma 1. *Given an undirected graph $G = (V, E)$ with m edges and n vertices and given a non-empty subset S of V , $p_S(u)$ and $d_G(u, p_S(u))$ can be computed in $O(m + n \log n)$ time over all vertices $u \in V$.*

Proof. Connect a new vertex s with a zero-weight edge to each vertex of S . Run Dijkstra (implemented with Fibonacci heaps) from s in this augmented graph. For each vertex $u \in V$, $p_S(u)$ is the unique ancestor of u belonging to S in the shortest path tree found and the distance from s to $p_S(u)$ in the tree equals $d_G(u, p_S(u))$. \square

Lemma 2. *Given an undirected graph $G = (V, E)$ with m edges and n vertices and given a non-empty subset S of V , G_S can be computed in $O(m + n \log n)$ time.*

Proof. Apply Lemma 1 to identify $p_S(u)$ and $d_G(u, p_S(u))$ for each $u \in V$. Then $E_S(u)$ can be found in time proportional to the degree of u . Hence, G_S can be found in $O(m)$ time in addition to the $O(m + n \log n)$ time from Lemma 1. \square

The following result is due to Baswana and Kavitha (see Lemmas 2.2 and 2.3 in [2]).

Lemma 3. *Let $G = (V, E)$ be an undirected n -vertex graph with positive edge weights and let $S \subseteq V$, $S \neq \emptyset$. For any two vertices $u, v \in V$, if $v \in B_S(u)$ then $d_{G_S}(u, v) = d_G(u, v)$. If S is obtained by picking each vertex independently with probability p , then E_S has expected size $O(n/p)$.*

3 A $(2k - 1)$ -Approximate Distance Oracle

In this section, we present a $(2k - 1)$ -approximate distance oracle with subquadratic preprocessing time for $k \geq 6$. As a warm-up, we first present an $O(k)$ stretch oracle with near-linear preprocessing. It is a trivial combination of the linear time spanner of Baswana and Sen [4] and the Thorup-Zwick oracle [15]. This idea is probably quite common knowledge and was noted by Sen [14]. However, the bounds obtained do not appear to be stated explicitly in the literature so we include them here. Later in this section and in Section 4, we shall refine this idea in order to get optimal tradeoff between size and stretch.

Theorem 1. *Let G be an undirected graph with m edges and n vertices and let $\epsilon > 0$ be a constant. For any integer $k \geq 1$, an $O(k)$ -approximate distance oracle for G of size $O(kn^{1+1/k})$ can be constructed in $O(km + kn^{1+1/k+\epsilon})$ time and can answer distance queries in $O(k)$ time.*

Proof. We compute in $O(km)$ time a spanner H of G with stretch at most $\lceil 1/\epsilon \rceil = O(1)$ and with $m_H = O(n^{1+1/\lceil 1/\epsilon \rceil}) = O(n^{1+\epsilon})$ edges using the linear time algorithm of Baswana and Sen [4]. We then construct a Thorup-Zwick oracle of stretch $2k - 1$ on top of H . It has size $O(kn^{1+1/k})$, query time $O(k)$, and preprocessing time $O(km_H n^{1/k}) = O(kn^{1+1/k+\epsilon})$. Since it has stretch $O(k)$ and H has stretch $O(1)$, the theorem follows. \square

3.1 Preprocessing

We now present our $(2k - 1)$ -approximate distance oracle and start with the preprocessing step. Each vertex is sampled with probability $p = n^{-i/k}$, for some $0 < i \leq k$ to be specified; we allow i to be a non-integer. Let S be the set of sampled vertices. We construct G_S in $O(m + n \log n)$ time using Lemma 2. The expected size of S is $pn = n^{1-i/k}$ and by Lemma 3, the expected size of E_S is $O(n/p) = O(n^{1+i/k})$. We can rerun the sampling until $|S| = \Theta(n^{1-i/k})$ and $|E_S| = O(n^{1+i/k})$; by Markov's inequality, only a constant expected number of reruns is needed for this.

We compute and store both $p_S(u)$ and $d_G(u, p_S(u))$ for all $u \in V$. By Lemma 1, this can be done in $O(m + n \log n)$ time.

Next, we build the distance oracle of Thorup and Zwick for the graph $G_S = (V, E_S)$. This takes $O(k|E_S|n^{1/k}) = O(kn^{1+(i+1)/k})$ expected time and requires $O(kn^{1+1/k})$ space. For some integer k' (to be specified), we apply the $O(k'm)$ time algorithm of Baswana and Sen [4] to find a $(2k' - 1)$ -spanner $H = (V, E_H)$ of G with $|E_H| = O(k'n^{1+1/k'})$ edges. For each pair of sampled vertices $p, q \in S$, we compute and store $d_H(p, q)$. The latter is done by running Dijkstra in H from each sampled vertex. Implementing Dijkstra with Fibonacci heaps, this takes a total of $O(|S|(|E_H| + n \log n)) = O(k'n^{2+1/k'-i/k})$ time. The space required to store all the $S \times S$ distances is $O(|S|^2) = O(n^{2-2i/k})$.

It follows from the above that total expected preprocessing time is $O(k'm + kn^{1+(i+1)/k} + k'n^{2+1/k'-i/k})$ and the amount of space needed for our oracle is $O(kn^{1+1/k} + n^{2-2i/k})$.

3.2 Answering a distance query

Now, let us consider how to answer a distance query for vertices $u, v \in V$, given the above preprocessing. In constant time, we look up vertices $p_S(u)$ and $p_S(v)$ as well as distances $d_G(u, p_S(u))$ and $d_G(v, p_S(v))$. We first query the distance oracle associated with G_S and get a distance estimate $\tilde{d}_1(u, v)$ in $O(k)$ time. We then obtain the precomputed value

$d_H(p_S(u), p_S(v))$ in constant time and obtain another distance estimate $\tilde{d}_2(u, v) = d_G(u, p_S(u)) + d_H(p_S(u), p_S(v)) + d_G(v, p_S(v))$. The smallest of $\tilde{d}_1(u, v)$ and $\tilde{d}_2(u, v)$ is then output as the answer to the query.

3.3 Bounding stretch

Let $\tilde{d}_G(u, v) = \min\{\tilde{d}_1(u, v), \tilde{d}_2(u, v)\}$ denote the distance estimate that our query step outputs. We show in the following that for a suitable choice of k' , $d_G(u, v) \leq \tilde{d}_G(u, v) \leq (2k - 1)d_G(u, v)$. The first inequality is clear since both $\tilde{d}_1(u, v)$ and $\tilde{d}_2(u, v)$ are the weights of some u - v -paths in G ; in particular, they are both at least as long as a shortest u - v -path. In the following, we show $\tilde{d}_G(u, v) \leq (2k - 1)d_G(u, v)$.

If $u \in B_S(v)$ or $v \in B_S(u)$ then by Lemma 3, $d_{G_S}(u, v) = d_G(u, v)$. In this case, the oracle for G_S outputs $\tilde{d}_1(u, v) \leq (2k - 1)d_G(u, v)$ so our query algorithm will output $\tilde{d}_G(u, v) \leq \tilde{d}_1(u, v) \leq (2k - 1)d_G(u, v)$, as desired.

Now assume that $u \notin B_S(v)$ and $v \notin B_S(u)$. Then $d_G(u, p_S(u)), d_G(v, p_S(v)) \leq d_G(u, v)$ and hence

$$\begin{aligned} d_H(p_S(u), p_S(v)) &\leq (2k' - 1)d_G(p_S(u), p_S(v)) \\ &\leq (2k' - 1)(d_G(p_S(u), u) + d_G(u, v) + d_G(v, p_S(v))) \\ &\leq (6k' - 3)d_G(u, v), \end{aligned}$$

so $\tilde{d}_2(u, v) = d_G(u, p_S(u)) + d_H(p_S(u), p_S(v)) + d_G(v, p_S(v)) \leq (6k' - 1)d_G(u, v)$. Hence, if we choose $k' = \lfloor k/3 \rfloor$, our query algorithm gives the desired stretch, i.e., $\tilde{d}_G(u, v) \leq \tilde{d}_2(u, v) \leq (2k - 1)d_G(u, v)$. Here we need to assume that $k \geq 3$ since k' needs to be at least 1.

3.4 Running time

Now, let us bound running time. Our query step runs in $O(k)$ time as mentioned above. With our choice of k' , preprocessing takes $O(km + kn^{1+(i+1)/k} + n^{2+1/\lfloor k/3 \rfloor - i/k})$ time. We consider three cases, $k \bmod 3 = 0$, $k \bmod 3 = 1$, and $k \bmod 3 = 2$, and we will fix a value of i in each of them that minimizes preprocessing time. In the first case, $k' = k/3$ and we set $i = k/2 + 1$, in the second case $k' = (k - 1)/3$ and we set $i = (k - 1)/2 + 3k/(2(k - 1))$, and in the third case $k' = (k - 2)/3$ and we set $i = (k - 2)/2 + (2k - 1)/(k - 2)$. The time and space bounds we obtain from these choices are stated in the following theorem.

Theorem 2. *Let G be an undirected graph with m edges and n vertices and let $k \geq 3$ be an integer. If $k \bmod 3 = 0$, a $(2k - 1)$ -approximate distance oracle for G of size $O(kn^{1+1/k})$ can be constructed in $O(km + kn^{3/2+2/k})$ time and can answer queries in $O(k)$ time. If $k \bmod 3 = 1$ resp. $k \bmod 3 = 2$, the same size and query bounds hold and construction time is $O(km + kn^{3/2+1/(2k)+3/(2(k-1))})$ resp. $O(km + kn^{3/2+2/(k-2)-1/(k(k-2))})$.*

We see that Theorem 2 breaks the quadratic time bound of Baswana and Kavitha [2] when $k \geq 6$.

4 Near-Linear Time Oracle

In this section, we give our $(2k - 1)$ -approximate distance oracle that breaks the preprocessing time of the previous section for larger k . It achieves a time bound arbitrarily close to linear

when k is sufficiently large. We shall modify the oracle of Section 3 and the modification we make is in the distance computations in the spanner H between sampled vertices in S . Instead of Dijkstra, we build an approximate distance oracle on top of H and then use it to report distances. However, we cannot apply this idea directly since we aim for $O(kn^{1+1/k})$ space; if we built the Thorup-Zwick oracle of that size on top of H , it would give stretch $2k - 1$ which for our oracle would be multiplied with the stretch of H . Instead, we will use the fact that the oracle only needs to report distances between vertices of S , allowing us to improve the Thorup-Zwick space/stretch tradeoff.

Overview of the Thorup-Zwick oracle: First, let us very briefly go through the algorithm of Thorup and Zwick. For an integer $\kappa \geq 1$, to build a size $O(\kappa n^{1+1/\kappa})$ oracle with stretch $2\kappa - 1$, sets A_0, \dots, A_κ are formed, with $V = A_0 \supseteq A_1 \supseteq A_2 \dots \supseteq A_\kappa = \emptyset$. For $i = 1, \dots, \kappa - 1$, set A_i is formed by picking each element of A_{i-1} independently with probability $n^{-1/\kappa}$. For $i = 0, \dots, \kappa - 1$, distance $d_G(A_i, v) = \min_{w \in A_i} d_G(w, v)$ is computed for each vertex v and a vertex $p_i(v)$ of A_i achieving this distance is kept. Then so called *bunches* $B(v)$ are formed around each vertex v and are defined by:

$$B(v) = \cup_{i=0}^{\kappa-1} \{w \in A_i \setminus A_{i+1} \mid d_G(w, v) < d_G(A_{i+1}, v)\}.$$

The oracle answers a query for the approximate distance between vertices u and v by repeatedly checking whether suitably chosen vertices belong to bunch $B(u)$ or bunch $B(v)$. More precisely, it starts by setting $w := u$ and initializes a counter $i := 0$. Then as long as $w \notin B(v)$, it increments i , swaps u and v and updates $w := p_i(u)$. When $w \in B(v)$ (which will happen at some point) $d_G(w, u) + d_G(w, v)$ is returned as approximate distance.

Improving space: The space requirement of the Thorup-Zwick oracle is dominated by the size of the bunches $B(v)$. For our application, we consider a *restricted* version of this oracle where we only care about queries between pairs of vertices from set S . It follows from the query step of the Thorup-Zwick oracle that we only need to store bunches $B(v)$ for $v \in S$ to answer such queries. As shown by Thorup and Zwick, each bunch has size $O(\kappa n^{1/\kappa})$. Hence, to build our restricted oracle, we only need $O(\kappa |S| n^{1/\kappa})$ space.¹

Our oracle: We get our near-linear time oracle by combining ideas from Section 3 with the restricted oracle above. As before, set S is formed by sampling each vertex with probability $p = n^{-i/k}$ and we construct the Thorup-Zwick oracle for graph $G_S = (V, E_S)$. We build a $(2k' - 1)$ -spanner $H = (V, E_H)$ with $O(k' n^{1+1/k'})$ edges. Now, instead of applying Dijkstra to find the exact distance $d_H(p, q)$ in H between each pair of samples $p, q \in S$, we build a restricted $(2\kappa - 1)$ -stretch oracle for H w.r.t. set S . It is built with the same procedure as that of Thorup and Zwick but only bunches $B(v)$ for $v \in S$ are kept. Total preprocessing for the oracles for G_S and H is then $O(k'm + kn^{1+(i+1)/k} + \kappa k' n^{1+1/k'+1/\kappa})$. We shall specify i , k' , and κ below.

To answer a distance query for vertices u and v , we do as in Section 3 but instead of using the exact spanner distance $d_H(p_S(u), p_S(v))$ we use that obtained by the restricted oracle.

¹In fact, using the source-restricted distance oracle in [13], we can reduce space further to $O(\kappa |S|^{1+1/\kappa})$ by keeping $O(|S|)$ bunches each of size $O(\kappa |S|^{1/\kappa})$. This also gives a slightly improved preprocessing time. Unfortunately, $|S|$ is too big for this result to give any significant improvement of our oracle.

The analysis for bounding the stretch for such a query is almost identical to that in Section 3.3. We only need to consider the case $u \notin B_S(v)$ and $v \notin B_S(u)$ and we get a stretch bounded by $2 + 3(2k' - 1)(2\kappa - 1)$ since we apply a $(2\kappa - 1)$ -approximate distance oracle instead of Dijkstra.

The size of our new oracle is bounded by the size of the oracle for G_S and the restricted oracle for H (note that we do not need to store H after the restricted oracle has been built). The first oracle requires $O(kn^{1+1/k})$ space and as we saw above, the restricted oracle requires $O(|S|\kappa n^{1/\kappa}) = O(\kappa n^{1-i/k+1/\kappa})$ space. Since we only allow space $O(kn^{1+1/k})$, we need $\kappa \geq k/(i+1)$. To minimize stretch, we pick κ as small as possible while satisfying this inequality, i.e., we pick $\kappa = \lceil k/(i+1) \rceil$. Substituting this for κ and requiring that stretch is at most $2k - 1$, we get the inequality

$$2 + 3(2k' - 1) \left(2 \left\lceil \frac{k}{i+1} \right\rceil - 1 \right) \leq 2k - 1 \Leftrightarrow k' \leq \frac{k + 3 \left(\left\lceil \frac{k}{i+1} \right\rceil - 1 \right)}{6 \left\lceil \frac{k}{i+1} \right\rceil - 3}.$$

To minimize running time, we pick k' as large as possible while satisfying this inequality, i.e., we set $k' = \lfloor \frac{k + 3(\lceil \frac{k}{i+1} \rceil - 1)}{6 \lceil \frac{k}{i+1} \rceil - 3} \rfloor$ (for this to make sense, we need $k' \geq 1$ but we shall choose i and k so that this is ensured).

We are aiming for a preprocessing time of $O(k'm + kn^{1+c/\sqrt{n}})$ for some constant $c > 0$. We shall pick c such that an i (possibly non-integer) can be found satisfying

$$\frac{2\sqrt{k}}{c} \leq \frac{k}{i+1} \leq \left\lceil \frac{k}{i+1} \right\rceil \leq \frac{c\sqrt{k}}{18}. \quad (1)$$

To ensure this, it suffices to require that $2\sqrt{k}/c + 1 \leq c\sqrt{k}/18$. Multiplying by c on both sides, we get

$$\frac{\sqrt{k}}{18}c^2 - c - 2\sqrt{k} \geq 0 \Leftrightarrow c \geq \frac{9}{\sqrt{k}} + 9\sqrt{\frac{1}{k} + \frac{4}{9}}.$$

Hence, since $k \geq 1$, picking $c = 9 + 3\sqrt{13}$ will allow us to pick an i satisfying inequalities (1). We claim that this choice of i gives the desired preprocessing time bound. Assume that $c/\sqrt{k} \leq 1$ since we are only interested in subquadratic bounds. Then

$$k' = \left\lfloor \frac{k + 3 \left(\left\lceil \frac{k}{i+1} \right\rceil - 1 \right)}{6 \left\lceil \frac{k}{i+1} \right\rceil - 3} \right\rfloor \geq \frac{k + 3 \left(\left\lceil \frac{k}{i+1} \right\rceil - 1 \right)}{\frac{c\sqrt{k}}{3} - 3} - 1 > \frac{k}{\frac{c\sqrt{k}}{3}} - 1 = \frac{3\sqrt{k} - c}{c} \geq \frac{2\sqrt{k}}{c}.$$

Since also $\kappa = \lceil k/(i+1) \rceil \geq k/(i+1) \geq 2\sqrt{k}/c$ and $\kappa k' = O(k)$, total preprocessing time is

$$O(k'm + kn^{1+(i+1)/k} + \kappa k' n^{1+1/k'+1/\kappa}) = O(\sqrt{k}m + kn^{1+c/(2\sqrt{k})} + kn^{1+c/\sqrt{k}}) = O(\sqrt{k}m + kn^{1+c/\sqrt{k}}),$$

as desired. We can now state the main result of the paper.

Theorem 3. *Let G be an undirected graph with m edges and n vertices and let $k \geq 1$ be an integer. Then a $(2k - 1)$ -approximate distance oracle for G of size $O(kn^{1+1/k})$ can be constructed in $O(\sqrt{k}m + kn^{1+c/\sqrt{k}})$ time, for some constant c , and can answer distance queries in $O(k)$ time.*

The bound we gave on constant c above was not very tight. In the following, let us bound this constant for large k . Instead of picking i such that inequalities (1) are satisfied, pick it such that $\kappa = \lceil k/(i+1) \rceil \geq k/(i+1) = \sqrt{k/6}$. Then

$$k' = \left\lfloor \frac{k + 3 \left(\left\lceil \frac{k}{i+1} \right\rceil - 1 \right)}{6 \left\lceil \frac{k}{i+1} \right\rceil - 3} \right\rfloor \geq \left\lfloor \frac{k + 3\sqrt{\frac{k}{6}} - 3}{6 \left\lceil \sqrt{\frac{k}{6}} \right\rceil - 3} \right\rfloor \geq \frac{k + 3\sqrt{\frac{k}{6}} - 3}{6\sqrt{\frac{k}{6}} + 3} - 1 = \frac{k - 3\sqrt{\frac{k}{6}} - 6}{\sqrt{6k} + 3},$$

which is at least 1 when $k \geq 29$. Hence

$$\frac{1}{\kappa} + \frac{1}{k'} \leq \sqrt{\frac{6}{k}} + \frac{\sqrt{6k} + 3}{k - 3\sqrt{\frac{k}{6}} - 6}$$

and it follows that we can pick c arbitrarily close to $2\sqrt{6}$ if k is bounded from below by a sufficiently large constant.

5 Concluding Remarks

For an undirected graph G with m edges, n vertices, and non-negative edge weights, and for an integer $k \geq 1$, the main result of this paper is a $(2k-1)$ -approximate distance oracle of G having size $O(kn^{1+1/k})$ and $O(k)$ query time which can be constructed in $O(\sqrt{k}m + kn^{1+c/\sqrt{k}})$ time for some constant c . We also gave an oracle with faster preprocessing for smaller k . Together, these two results break the quadratic preprocessing time of Baswana and Kavitha for $k \geq 6$ and improve the $O(kmn^{1/k})$ bound of Thorup and Zwick when G is not too sparse and k not too small.

Assuming the girth conjecture, our oracle is optimal in every respect when $m = \Omega(n^{1+c/\sqrt{k}})$ and $k = O(1)$ since then preprocessing time is linear. Whether linear preprocessing time is achievable also for sparser graphs remains an open problem. Our oracles break the quadratic preprocessing bound when $k \geq 6$. What is possible for smaller k ? Finally, we believe that by using existing techniques, it should be possible to derandomize our oracles at only a small increase in preprocessing time.

References

- [1] S. Baswana, A. Gaur, S. Sen, and J. Upadhyay. Distance oracles for unweighted graphs: Breaking the quadratic barrier with constant additive error. In Proc. 35th International Colloquium on Automata, Languages and Programming (ICALP), pp. 609–621, 2008.
- [2] S. Baswana and T. Kavitha. Faster Algorithms for All-Pairs Approximate Shortest Paths in Undirected Graphs. SIAM J. Comput., Vol. 39, No. 7, pp. 2865–2896, 2010.
- [3] S. Baswana, T. Kavitha, K. Mehlhorn, and S. Pettie. New constructions of (α, β) -spanners and purely additive spanners. In Proc. 16th ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 672–681, 2005.
- [4] S. Baswana and S. Sen. A Simple and Linear Time Randomized Algorithm for Computing Sparse Spanners in Weighted Graphs. Random Structures & Algorithms, 30 (2007), pp. 532–563.

- [5] T. M. Chan. More algorithms for all-pairs shortest paths in weighted graphs. In Proc. 39th Annual ACM Symposium on Theory of Computing (STOC), pp. 590–598, 2007.
- [6] E. Cohen and U. Zwick. All-pairs small stretch paths. *J. Algorithms*, 38 (2001), pp. 335–353.
- [7] D. Dor, S. Halperin, and U. Zwick. All-pairs almost shortest paths. *SIAM J. Comput.*, 29(5):1740–1759, 2000.
- [8] M. Elkin and D. Peleg. $(1 + \epsilon, \beta)$ -spanner constructions for general graphs. *SIAM J. Comput.*, 33(3):608–631, 2004.
- [9] P. Erdős. Extremal problems in graph theory. In *Theory of Graphs and its Applications* (Proc. Sympos. Smolenice, 1963), Czechoslovak Acad. Sci., Prague, 1964, pp. 29–36.
- [10] M. Mendel and A. Naor. Ramsey partitions and proximity data structures. *Journal of the European Mathematical Society*, 9(2):253–275, 2007. See also FOCS’06.
- [11] M. Pătraşcu and L. Roditty. Distance Oracles Beyond the Thorup-Zwick Bound. In Proc. 51st Annual IEEE Symposium on Foundations of Computer Science (FOCS), pp. 815–823, 2010.
- [12] S. Pettie. A new approach to all-pairs shortest paths on real-weighted graphs. *Theoret. Comput. Sci.*, 312 (2004), pp. 47–74.
- [13] L. Roditty, M. Thorup, and U. Zwick. Deterministic constructions of approximate distance oracles and spanners. L. Caires et al. (Eds.): *ICALP 2005*, LNCS 3580, pp. 661–672, 2005.
- [14] S. Sen. Approximating Shortest Paths in Graphs. S. Das and R. Uehara (Eds.): *WALCOM 2009*, LNCS 5431, pp. 32–43, 2009.
- [15] M. Thorup and U. Zwick. Approximate Distance Oracles. *J. Assoc. Comput. Mach.*, 52 (2005), pp. 1–24.
- [16] M. Thorup and U. Zwick. Spanners and emulators with sublinear distance errors. In Proc. 17th ACM-SIAM Symposium on Discrete Algorithms (SODA), pp. 802–809, 2006.
- [17] U. Zwick. All-pairs shortest paths using bridging sets and rectangular matrix multiplication. *J. Assoc. Comput. Mach.*, 49 (2002), pp. 289–317.